

# Stanford Machine Learning Week 2. § Linear Regression with multiple variables.

- multiple inputs (information) to predict.
- input:  $x_1, x_2, \dots, x_n$ ; output =  $y$ .

e.g.

$x_1$ size	$x_2$ # of bedrooms	$x_3$ Number of floors	$x_4$ Age of home	price ( $y$ )
				m samples

Notation:

$n$  - number of features.

$x^{(i)}$  - input (features) of  $i$ th training sample

$x_j^{(i)}$  - value of feature  $j$  in  $i$ th training sample

$$\Rightarrow x^{(i)} = [x_j^{(i)}].$$

index to training set

$$x^{(i)} = \begin{pmatrix} x^{(i)}_1 \\ x^{(i)}_2 \\ \vdots \\ x^{(i)}_n \end{pmatrix}$$

$n$ -dimensional vector.

## - Hypothesis

previously  $h_\theta(x) = \theta_0 + \theta_1 x$

now  $\rightarrow h_\theta(x) = \theta_0 + \theta_1 x_1 + \dots + \theta_n x_n$

define:  $x_0 = 1$ , i.e.  $x_0^{(i)} = 1 \quad \forall i$

$$\Rightarrow x = \begin{bmatrix} x_0 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}; \quad \theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^n.$$

$$h_\theta(x) = \theta_0 x_0 + \dots + \theta_n x_n$$

$$= \theta^T \cdot x \quad (\text{expressed as inner product of vectors})$$

$$\Rightarrow \text{Multivariate Linear Regression: } h_\theta(x) = \theta^T \cdot x.$$

# Linear Regression with Multiple Variables - Gradient descent for Multiple Variables

Recall: Hypothesis  $h_{\theta}(x) = \underline{\theta}^T \cdot \underline{x} = \sum \theta_i x_i$

Parameters  $\underline{\theta} = \begin{pmatrix} \theta_0 \\ \vdots \\ \theta_n \end{pmatrix} \rightarrow \in \mathbb{R}^{n+1}.$

Cost function.  $J(\underline{\theta}) = \frac{1}{2m} \sum_{i=1}^m [h_{\theta}(\underline{x}^{(i)}) - y^{(i)}]^2$

Gradient descent for multiple variables.

1) repeat {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\underline{\theta}) \quad [\text{simultaneously update } \forall j=0 \dots n]$$

}

=> 2) repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m [h_{\theta}(\underline{x}^{(i)}) - y^{(i)}] \cdot x_j^{(i)}$$

}

$\nearrow$  i denotes sample #

$\hookrightarrow$  j denotes features.

Note:  $x_0^{(i)} = 1$ . by def'n.

# Linear Regression with multiple variables - Gradient descent

## Feature Scaling

- Idea: make sure features are on similar schedule.
- Get every feature into  $-1 \leq x_i \leq 1$  range. (wonder)

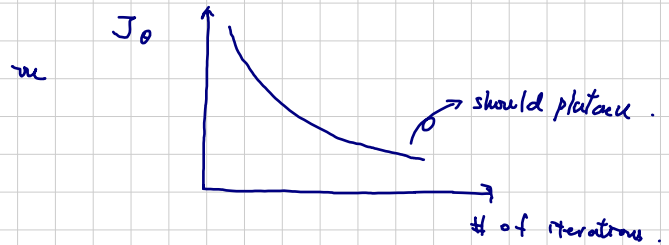
## Mean Normalization

- $x_i \rightarrow x_i - \mu_i$  ( $E(x_i - \mu_i) = 0$ ,  $x_i - \mu_i$  has zero mean)
- standard normal:  $x_i \rightarrow \frac{x_i - \mu_i}{s_i}$   
↳ approximate standard deviation by range. (max-min)

## Learning Rate

- make sure gradient descent is working  $\rightarrow$  plot.  $\min_{\theta} J_{\theta}$  over # of iterations.

- for sufficiently small  $\alpha$ ,  
 $J_{\theta}$  should decrease after every iteration



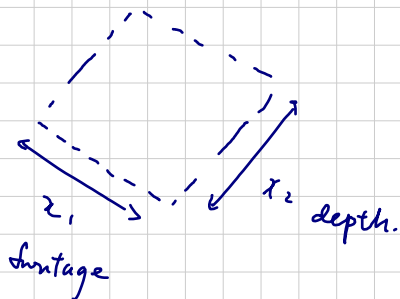
## Summary

- $\alpha$ 
  - $\rightarrow$  too small: slow convergence
  - $\rightarrow$  too large: may not converge

- to choose  $\alpha$ , try 0.001, 0.01, 0.1, 1.

## Features and Polynomial Regression

ex. housing price prediction.



$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 \quad ; \quad \text{alternatively let } x = \text{frontage} \times \text{depth (area)}$$
$$\text{let } h_{\theta}(x) = \theta_0 + \theta_1 \cdot x$$

• Polynomial regression: "fitting" into different polynomials. by choice

- We can approach this using multivariate regression.

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$
$$\begin{array}{ccc} \downarrow & \downarrow & \downarrow \\ x_1 = x & x_2 = x^2 & x_3 = x^3 \end{array}$$

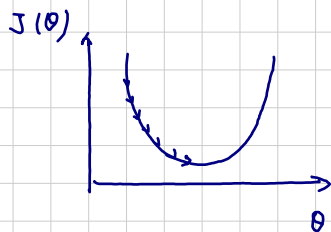
- note: apply "feature scaling" to each of  $x_1, x_2, x_3$  !

{ Computing Parameters Analytically

### Normal Equation

• Gradient descent.

↳ iterative approach



• to find  $\min J$ .

↳ set  $\frac{\partial}{\partial \theta_j} J(\theta) = 0$  ; solve for  $\theta_0, \theta_1, \dots, \theta_n$ .

• Normal method: solving for  $\theta$  analytically.

• construct matrix.  $X = \begin{pmatrix} x_0 & x_1 & \dots & x_m \end{pmatrix} \begin{array}{l} \text{ } \\ \text{ } \\ \text{ } \end{array} \left. \vphantom{\begin{pmatrix} x_0 & x_1 & \dots & x_m \end{pmatrix}} \right\} \begin{array}{l} n+1; n \text{ features.} \\ m \text{ samples} \end{array}$

$y = \begin{pmatrix} \vdots \end{pmatrix}$   $m$  dimensional vector

$$\theta = (X^T X)^{-1} \cdot X^T y$$

• General case:  $m$  examples  $(x^{(1)}, y^{(1)})$ , ...,  $(x^{(m)}, y^{(m)})$ ;  $n$  features

$$x^{(i)} = \begin{pmatrix} x_0^{(i)} \\ \vdots \\ x_n^{(i)} \end{pmatrix} \in \mathbb{R}^{n+1}; \text{ Design matrix } X = \begin{bmatrix} - [x^{(1)}]^T - \\ \vdots \\ - [x^{(m)}]^T - \end{bmatrix} \in \mathbb{R}^{m \times (n+1)}$$

e.g. if  $x^{(i)} = \begin{pmatrix} 1 \\ x_1^{(i)} \end{pmatrix}$ ,  $X = \begin{pmatrix} 1 & x_1^{(1)} \\ \vdots & \vdots \\ 1 & x_1^{(m)} \end{pmatrix}$ ;  $y = \begin{pmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{pmatrix}$

$$\theta = (X^T X)^{-1} X^T y$$

$(X^T X)^{-1}$  is inverse of  $X^T X$

• Gradient descent

Normal Equation

- need to choose  $\alpha$
- need many iterations
- works well w/ large  $n$ .  
(10<sup>6</sup>)

- need to compute  $(X^T X)^{-1}$
- ↓
- slow if  $n$  is large. ( $O(n^3)$  for inverse computation)

### Normal equation and non-invertability

Normal equation:  $\theta = (X^T X)^{-1} X^T y$

- What if  $X^T X$  is not-invertable.

- (1) Redundant features (linearly dependent)
- (2) Too many features ( $m \leq n$ ), more than sets.  
↳ delete some features, or use regularization. (later)

